

# Naïve Bayes classifier

A small weather data set on previous records of

- (i) weather conditions and
- (ii) whether certain event happens (i.e. certain activity is “played”)

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

## A new case for prediction:

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

Play=?

- **E**: The evidences (observations) we have:

$$E_1 = \frac{\text{Outlook}}{\text{Sunny}}, E_2 = \frac{\text{Temp.}}{\text{Cool}}, E_3 = \frac{\text{Humidity}}{\text{High}}, E_4 = \frac{\text{Windy}}{\text{True}}$$

$$\mathbf{E} = (E_1, E_2, E_3, E_4) =$$

Outlook	Temp.	Humidity	Windy
Sunny	Cool	High	True

- **H**: whether the event happens:  
Two possible predictions (hypotheses):

**H** = “Play=yes”    or

**H** = “Play=no”

## The basic probabilistic approach:

Compare two conditional probabilities

$\Pr( H = \text{“Play=yes”} \mid E =$ 

Outlook	Temp.	Humidity	Windy
Sunny	Cool	High	True

 $)$  versus

$\Pr( H = \text{“Play=no”} \mid E =$ 

Outlook	Temp.	Humidity	Windy
Sunny	Cool	High	True

 $)$  .

It is hard to directly estimate these two probabilities from the data set.  
(Why?)

## The Naïve Bayes approach:

(I) Use the Bayes rule

- Probability of event  $H$  given evidence  $E$ :

$$Pr[H|E] = \frac{Pr[E|H]Pr[H]}{Pr[E]}$$

- *A priori* probability of  $H$  :  $Pr[H]$ 
  - Probability of event *before* evidence is seen
- *A posteriori* probability of  $H$  :  $Pr[H|E]$ 
  - Probability of event *after* evidence is seen

**Thomas Bayes**

**Born:** 1702 in London, England

**Died:** 1761 in Tunbridge Wells, Kent, England



## The Naïve Bayes approach:

Just calculate  $Pr[E|H]Pr[H]$  as the likelihood and compare

$$Pr(E = \begin{array}{|c|c|c|c|} \hline \text{Outlook} & \text{Temp.} & \text{Humidity} & \text{Windy} \\ \hline \text{Sunny} & \text{Cool} & \text{High} & \text{True} \\ \hline \end{array} \mid H = \text{“Play=yes”}) *$$

$$Pr(H = \text{“Play=yes”})$$

with

$$Pr(E = \begin{array}{|c|c|c|c|} \hline \text{Outlook} & \text{Temp.} & \text{Humidity} & \text{Windy} \\ \hline \text{Sunny} & \text{Cool} & \text{High} & \text{True} \\ \hline \end{array} \mid H = \text{“Play=no”}) *$$

$$Pr(H = \text{“Play=no”}).$$

No need to worry about  $Pr[E]$  since

$Pr(E = \begin{array}{|c|c|c|c|} \hline \text{Outlook} & \text{Temp.} & \text{Humidity} & \text{Windy} \\ \hline \text{Sunny} & \text{Cool} & \text{High} & \text{True} \\ \hline \end{array})$  is the same denominator on the right hand side of the Bayes rule.

## The Naïve Bayes approach:

- (II) Use the Naïve assumption on independency: Individual evidences ( $E_1, E_2, E_3, E_4 \dots$ ) are independently affected by the underlying event separately.

$$\Pr[E|H]$$

$$= \Pr[E_1|H] * \Pr[E_2|H] * \dots * \Pr[E_n|H]$$

It is much easier estimate conditional probabilities

$\Pr(E_1 | H)$ ,  $\Pr(E_2 | H)$ ,  $\Pr(E_3 | H)$ , and  $\Pr(E_4 | H)$  from the data set . (Why?)

How to estimate  $\Pr(E_1 | H)$ ,  $\Pr(E_2 | H)$ ,  $\Pr(E_3 | H)$ ,  $\Pr(E_4 | H)$  ?

Outlook			Temperature			Humidity			Windy			Play	
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Estimate  $\Pr(E_1 | H)$  and  $\Pr(E_2 | H)$

Outlook			Temperature			Humidity			Windy		Play		
	Yes	No		Yes	No		Yes	No	Yes	No	Yes	No	
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/	5/
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5	14	14
Rainy	3/9	2/5	Cool	3/9	1/5								

$$\Pr(E_1 = \frac{\text{Outlook}}{\text{Sunny}} | H = \text{“Play=yes”}) : 2/9 \text{ (why?)}$$

$$\Pr(E_1 = \frac{\text{Outlook}}{\text{Sunny}} | H = \text{“Play=no”}) : 3/5 \text{ (why?)}$$

$$\Pr(E_2 = \frac{\text{Temp.}}{\text{Cool}} | H = \text{“Play=yes”}) : 3/9 \text{ (why?)}$$

$$\Pr(E_2 = \frac{\text{Temp.}}{\text{Cool}} | H = \text{“Play=no”}) : 1/5 \text{ (why?)}$$



Estimate  $\Pr(E_3 | H)$  and  $\Pr(E_4 | H)$

Outlook			Temperature			Humidity			Windy			Play	
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

$$\Pr(E_3 = \begin{array}{c} \text{Humidity} \\ \text{High} \end{array} | H = \text{“Play=yes”}) : 3/9 \text{ (why?)}$$

$$\Pr(E_3 = \begin{array}{c} \text{Humidity} \\ \text{High} \end{array} | H = \text{“Play=no”}) : 4/5 \text{ (why?)}$$

$$\Pr(E_4 = \begin{array}{c} \text{Windy} \\ \text{True} \end{array} | H = \text{“Play=yes”}) : 3/9 \text{ (why?)}$$

$$\Pr(E_4 = \begin{array}{c} \text{Windy} \\ \text{True} \end{array} | H = \text{“Play=no”}) : 3/5 \text{ (why?)}$$

Estimate Pr(H):

Outlook			Temperature			Humidity			Windy		Play		
	<i>Yes</i>	<i>No</i>		<i>Yes</i>	<i>No</i>		<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5	14/14	14/14
Rainy	3/9	2/5	Cool	3/9	1/5								

$\Pr(H = \text{"Play=yes"}) = 9/14$  (why?)

$\Pr(H = \text{"Play=no"}) = 5/14$  (why?)

Just calculate the likelihood  $Pr[E|H]Pr[H]$

for  $H = \text{“Play=yes”}$  and for  $H = \text{“Play=no”}$

For example,

$Pr(E =$ 

Outlook	Temp.	Humidity	Windy
Sunny	Cool	High	True

 $| H = \text{“Play=yes”}) *$

$Pr(H = \text{“Play=yes”})$

=

$$\begin{aligned} &Pr[Outlook = Sunny|yes] \\ &\times Pr[Temperature = Cool|yes] \\ &\times Pr[Humidity = High|yes] \\ &\times Pr[Windy = True|yes] \\ &\times \underline{Pr[yes]} \end{aligned}$$

The results:

Likelihood of the two classes

$$\text{For "yes"} = 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$\text{For "no"} = 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0206$$

Conversion into a probability by normalization:

$$P(\text{"yes"}) = 0.0053 / (0.0053 + 0.0206) = 0.205$$

$$P(\text{"no"}) = 0.0206 / (0.0053 + 0.0206) = 0.795$$

The prediction: "Play=no"